

並列構文解析法とその評価

(財) 新世代コンピュータ技術開発機構

松本裕治 (Yuji Matsumoto)

1. はじめに

本稿の目的は、現在開発中の自然言語処理用の並列構文解析システム A X (Analyser for natural language syntaX) の評価を行なうことである。このシステムは、Parlog [Clark 84]、Concurrent Prolog [Shapiro 83]、G H C (Guarded Horn Clauses) [Ueda 85] などの並列論理型言語上での実現を目指した構文解析システムである。このシステムは、構文上は、文脈自由文法を扱うが、論理型言語より派生した自然言語用の文法記述形式である D C G (Definite Clause Grammars) [Pereira 80] で書かれた文法を処理することができるため、実際には、文脈自由文法以上の文法をも取り扱うことが可能である。

本稿では、構文解析アルゴリズムそのものの評価を行なうため、文法を文脈自由文法に限って、その計算量の複雑さについて、考察する。

まず、最初に、アルゴリズムの記述を行ない、並列計算での複雑さを評価する。次に、アルゴリズムの並列論理型言語による実現法を示す。その結果、本手法は、入力文の長さに対して

線形時間の並列構文解析プログラムであり、しかも、同一計算を繰り返さない最適アルゴリズムであることが分かる。

また、本アルゴリズムは、逐次処理マシンにとっても効率よく、Prolog によって実用的な構文解析システムとして、利用することができる。アルゴリズム自体は、AX と呼ばれるが、Prolog で逐次マシン上に実現されたシステムを SAX (Sequential AX)、GHC 等の並列型言語によって実現されたシステムを PAX (Parallel AX) と呼ぶ。

2. アルゴリズムの記述

並列アルゴリズムの記述には、並列計算モデルが必要である。本節では、我々のアルゴリズムの記述に、W-RAM (Write-enabled Parallel Random Access Machine) を用いる。これは、自然数によって、インデクシングされた、共有記憶を持ち、各記憶セルに対して、同時読み出し、および、同時書き込みが許されている並列 RAM である。ただし、同時書き込みが行なわれる場合には、書き込まれる内容が同じでなければならない。

文脈自由文法を次のように定義しよう。 $G = (V_N, V_T, P, S)$ 。ここに、 V_N は非終端記号の有限集合、 V_T は終端記号の有限集合、 P は文法記号の有限集合、 S は初期記号を表わす非終端記号である。以下、入力文字列を、 $a_1 a_2 \dots a_n$ であらわす。また、一般に、非終端記号を大文字のアルファベットで、終端記号を小文字のアルファベットで、非終端記号と終端記号よりなる記号列をギリシャ文字で表わす。また、文法規則の右辺に現れる任意の連続する二つの記号の間に識別子として、

個別に自然数を割り当て、そのような自然数の集合を K とする。
 例えば、 $A \rightarrow B C D$ 、 $E \rightarrow F G$ を P に含まれる文法規則とすると、これには、 $A \rightarrow B_1 C_2 D$ 、 $E \rightarrow F_3 G$ のように識別用の自然数が割り当てられているとする。

準備として、次の二種類の集合を定義する。これらはアルゴリズムでは、表のような構造をもつ記憶領域として扱われる。一つは、

$$S = \{ (A, i, j) : A \in V_N, 0 \leq i < j \leq n \}$$

なる集合。もう一つは、 $1 \leq j \leq n - 1$ なる j に対して、

$$T_j = \{ (k, i) : k \in K, 0 \leq i < n \}$$

として定義される一連の集合である。

S に含まれる要素 (A, i, j) の意味は、非終端記号 A が、入力文字列の部分列 $a_{i+1} \dots a_j$ を葉 (leaves) とする解析木の根 (root) になり得ることである。 T_j 内の要素 (k, i) の意味は、 k が、文法規則 $A \rightarrow \alpha k \beta$ に含まれる識別子とすると、この規則の右辺で k によって区切られた箇所の直前までの記号列、すなわち α が、入力文字列の部分列 $a_{i+1} \dots a_j$ によって完成されたということである。これらの集合は各要素を項目として持つ表として取り扱われ、全ての要素は初期値 'false' を持つものとする。

さて、文脈自由文法についての並列構文解析のアルゴリズムを示そう。

アルゴリズム AX

begin

```

1: for each  $0 \leq i < n$ ,  $A \in V_N$  such that  $A \rightarrow a_{i+1}$ 
    parallel do  $(A, i, i+1) := \text{true}$  ;
repeat  $N$  times
2: for each  $(B, i, j) = \text{true}$  and  $A \rightarrow B \ k \ \alpha$ 
    parallel do  $(k, i) := \text{true}$  in  $T_j$  ;
2': for each  $(B, i, j) = \text{true}$  and  $A \rightarrow B$ 
    parallel do  $(A, i, j) := \text{true}$  ;
3: for each  $(B, j, m) = \text{true}$  and  $(k, i) = \text{true}$  in  $T_j$ 
    such that  $A \rightarrow \alpha \ k \ B$ 
    parallel do  $(A, i, m) := \text{true}$  ;
3': for each  $(B, j, m) = \text{true}$  and  $(k, i) = \text{true}$  in  $T_j$ 
    such that  $A \rightarrow \alpha \ k \ B \ k+1 \ \beta$ 
    parallel do  $(k+1, i) := \text{true}$  in  $T_m$ 
end repeat ;
if  $(S, 0, n) = \text{true}$  then ACCEPT
end ;

```

2 および 3, 3' によって、文法規則の右辺に含まれる記号が順次処理されていくので、文法規則の右辺に現れる記号数の最大値を c とすると、アルゴリズムの繰り返し回数 N は、 $(c-1)n$ 回でよく、このアルゴリズムは $o(n)$ の線形アルゴリズムとなる。本アルゴリズムは、チャートパーズング [Kay 80] を上昇型に並列に行なうのと同等で、各部分木に対して、文法規則の適用は、重複しては行なわれないから、同一処理の繰り返しは避けられている。

3. 並列論理型言語による実現

本節では、上で述べたアルゴリズムを並列論理型プログラミング言語で実現する方法を示す。記述言語としては、GHCを用いる。

並列論理型言語によるプログラミングは、ストリームを介して通信し合う並列プロセスの記述を基本とする。そこで、 (A, i, j) なる項をプロセスとして論理型言語では述語として定義し、 T_j をプロセス間を流れるストリームとして取り扱うことを考える。入力文字列 $a_1 a_2 \dots a_n$ に対して、初期実行を行なう呼び出しを

$$a_1(T_0, T_1), a_2(T_1, T_2), \dots, a_n(T_{n-1}, T_n).$$

と定義する。これは、GHCで `and` 並列に動くプロセスである。ここに、 T_0, T_n は、システムの使用者によって定義され、使用される、定数またはストリームで、その使い方についてはここでは省略する。

前節で示したアルゴリズムの各ステップは、次のようなGHC節に一对一に変換することができる。ただし、論理型言語の慣習により、述語名は英小文字で、変数は英大文字を先頭に持つ項として記述されている。

- 1: $a_i(T_{i-1}, T_i) :- a(T_{i-1}, T_i). \quad (\text{for each } A \rightarrow a_i \in P)$
- 2: $b(T_i, T_j) :- T_j = [(k, T_i)]. \quad (\text{for each } A \rightarrow B \text{ k } \alpha \in P)$
- 2': $b(T_i, T_j) :- a(T_i, T_j). \quad (\text{for each } A \rightarrow B \in P)$
- 3: $b([(k, T_i) | T_j], T_m) :- a(T_i, T_m'), b(T_j, T_m'')$
 $\quad \text{merge}(T_m', T_m'', T_m).$

(for each $A \rightarrow \alpha \quad k \quad B \in P$)

3': $b([(k, T_i) | T_j'], T_m) :- T_m = [(k+1, T_i) | T_m'], b(T_j', T_m')$.

(for each $A \rightarrow \alpha \quad k \quad B \quad k+1 \quad \beta \in P$)

このようにして定義されたそれぞれの G H C 節は、独立に動くことができ、同一名をもつ述語の出力（第 2 引数）は、マージされて、次のプロセスへ渡すことによってこれらを a n d 並列に実行することが可能である。これらの G H C 節の定義を見て分かるように、文法規則の右辺に現れる記号（終端記号、非終端記号を問わず）一つずつに対して、対応する G H C 節が生成される。同一名をもつ述語は、全て a n d 並列に動くプロセスとして定義できる。ただし、3 と 3' は、互いに排他的であるので、実際には、これらの節は、o r 並列に実行することも可能であり、現システムでは、そのように定義されている [松本 86]。具体的な文法規則について、G H C 節への変換の例を付録に示す。

4. 論理型言語による実現についての問題点

本節では、A X の実現に関して述べておかねばならない問題点について考察する。

同時書き込みについての問題は、論理型言語による実現では部分解析木がプロセスとして定義されているため、同一プロセスの生成の問題に置き換わっている。資源が無尽蔵と考えてよいならば、これは計算時間には影響を与えない。ただし、現実には、プロセッサ数は有限と考えねばならないので、その場合には、同一プロセスの生成を阻止するような、メタな機構を考

えねばならない。ただし、本システムが単なる文脈自由文法ではなく、DCGで書かれた自然言語文法を扱うことを考えると、構文解析だけではなく、それに伴う意味解析等も考慮に入れなければならない。従って、入力文の同一部分列からたとえ同じ非終端記号を根とする解析木が生成されようとも、文法規則に同一のものが含まれていない限り、それらの内容はどこかが違うのであるから、同じ非終端記号が構成されてしまっても構わないという考え方もある。この問題は、資源の有効利用と重複処理のオーバーヘッドとが絡み、今後解決しなければならない問題である。

また、いろいろなプロセスから生み出される出力が、ストリームによって、一本にまとめられてしまうため、並列性がある程度失われてしまうように見えるが、これは、ストリームを並列論理型言語の中でどのように実現するかには依存する問題であり、ストリームがランダムアクセス可能な記憶構造として実現されれば、効率に対する影響は軽減される。この問題に関してだけでなく、ストリームに関する処理の効率は並列論理型言語にとって極めて重要な問題であることを指摘しておきたい。

本アルゴリズムによる文の解析は、完全に上昇型 (bottom-up) に行なわれるが、下降型 (top-down) の戦略も部分的に利用することができる。実際のシステムでは、上からの予測を、ストリームの中から無駄な要素を取り除くフィルターのようなプロセスとして実現されている。

5. 結論と考察

文脈自由文法に対する並列構文解析アルゴリズムとして、A X は、線形時間のアルゴリズムであることが分かった。文脈自由文法の認識についての並列アルゴリズムとしては、他に、時間的により優れたものが提案されている。Alternating Turing Machine [Chandra 81] や P - R A M (Parallel Random Access Machine) [Fortune 78] では、 $o(\log^2 n)$ のアルゴリズムが報告されている [Ruzzo 80] [Rytter 84]。また、W - R A M では、それらのアルゴリズムを $o(\log n)$ でシミュレートできる。

しかし、これらのアルゴリズムでは、無駄なプロセスが相当作られるので、現実の構文解析システムとして使うには難がある。

一方、我々のアルゴリズムでは、無駄なプロセスやデータの生成は、極力おさえられている。逐次的実行を行なった場合でも、その処理の複雑さは、チャートパーザと同等もしくはそれ以下であり、データやプロセスの生成量は最適化されている。このシステムは、Prolog 上でもインプリメントされ、自然言語処理用の構文解析システムとして利用される予定である。

今後の課題としては、無駄な計算を抑えながら、如何に計算時間を短縮するかということと、同一処理を防ぐためのメタな機能をどのように自然に取り込むかを考えなければならない。

【参考文献】

[Chandra 81] Chandra, A.K., et al., "Alternation,"

JACM, vol.28, no.1, pp.114-133, 1981.

- [Clark 84] Clark, K. and Gregory, S., "PARLOG: Parallel Programming in Logic," Imperial College Research Report DOC 81/16, July 1984.
- [Kay 80] Kay, M., "Algorithm Schemata and Data Structures in Syntactic Processing," XEROX PARC, CSL-80-12, Oct. 1980.
- [Fortune 78] Fortune, S. and Wyllie, J., "Parallelism in Random Access Machines," Proc. 10th ACM Symposium on Theory of Computation, pp.114-118, 1978.
- [Pereira 80] Pereira, F. and Warren, D., "Definite Clause Grammars for Language Analysis," Artificial Intelligence 13, pp.231-278, 1980.
- [Ruzzo 80] Ruzzo, W.L., "Tree-size Bounded Alternation" Journal of computer and System Sciences 21, pp.218-235, 1980.
- [Rytter 85] Rytter, W., "Parallel Time $o(\log N)$ Recognition of Unambiguous CFLs," Lecture Notes in Computer Science 199, pp.380-389, 1985.
- [Shapiro 83] Shapiro, E., "A Subset of Concurrent Prolog and its Interpreter," Technical Report TR-003, ICOT, Feb.1983.
- [Ueda 85] Ueda, K., "Guarded Horn Clauses," Proc. The Logic Programming Conference '85, ICOT, pp.225-236, July 1985.

[松本 86] 松本裕治、「並列構文解析」、情報処理学会
自然言語処理研究会、1986年1月

【付録】

文脈自由文法から G H C 節への変換の例を示す。(1)で示すのは、文法規則の一部であり、(2)が、変換後のプログラムの内 noun に関する G H C 節だけを取り出したものである。ただし、文法規則中の右辺の非終端記号間に書かれている自然数は、識別子である。

(1) NP → DET 1 NOUN

NP → DET 2 NOUN 3 REL_CLAUSE

NP → NOUN

NOUN → NOUN 4 PP

NOUN → NOUN 5 PRED

(2) noun(X,Y) :- noun1(X,Y1),noun2(X,Y2),

merge(Y1,Y2,Y).

noun1(X,Y) :- Y=[(4,X),(5,X)|Y1],np(X,Y1).

noun2([],Y) :- Y=[].

noun2([(1,X)|X1],Y) :- np(X,Y1),noun2(X1,Y2),

merge(Y1,Y2,Y).

noun2([(2,X)|X1],Y) :- Y=[(3,X)|Y1],noun2(X1,Y1).

noun2([_IX],Y) :- otherwise | noun2(X,Y).